

1 This application is submitted in the name of the following inventor(s):

2

3	<i>Inventor</i>	<i>Citizenship</i>	<i>Residence City and State</i>
4	Scott SCHOENTHAL	United States	San Ramon, California

5

6 The assignee is *Network Appliance, Inc.*, a California corporation having an
7 office at 495 East Java Drive, Sunnyvale CA 94089.

8

9 TITLE OF THE INVENTION

10

11 Persistent and Reliable Delivery of Event Messages

12

13 BACKGROUND OF THE INVENTION

14

15 1. *Field of the Invention*

16

17 This invention relates to persistent and reliable delivery of event messages,
18 including event messages in file server systems in which it is desired to maintain reliable
19 file system consistency.

20

2. *Related Art*

In systems that provide services to clients, such as those including file servers and similar devices, it often occurs that the system, or some subsystem within that system, generates a message indicating the occurrence of a special event. Typically, the special event is an error of some kind, and the message conveys information regarding the nature of the special event, such as the type of error and the subsystem within which the error occurred. Many systems that provide services, including file servers, make efforts to assure that the services are reliably provided, and that the system providing the services is in a consistent state at all times. Thus, such systems find it advantageous to assure that all state information regarding the system, including state information relating to error messages, is persistently and reliably maintained. Such systems also find it advantageous to assure that all event messages are reliably delivered and persistently maintained until delivery is confirmed by the intended recipient of the event message.

Accordingly, it would be advantageous to provide a technique for persistent and reliable delivery of event messages, that is not subject to the drawbacks of the known art. Preferably, those parts of the system responsible for delivering event messages are able to persistently maintain those event messages until delivery of those event messages has been confirmed by the intended recipient of the event message. Moreover, those parts of the system responsible for recovering from system crashes and other system er-

1 rors are able to persistently maintain those event messages until delivery, even after re-
2 covery from system crashes or other system errors.

4 SUMMARY OF THE INVENTION

6 The invention provides a method and system for persistent and reliable de-
7 livery of event messages, that is not subject to the drawbacks of the known art. Those
8 parts of the system responsible for delivering event messages are able to persistently
9 maintain those event messages until delivery of those event messages has been confirmed
10 by the intended recipient of the event message. Those parts of the system responsible for
11 recovering from system crashes and other system errors are able to persistently maintain
12 those event messages until delivery, even after recovery from system crashes or other
13 system errors.

15 In a first aspect of the invention, the system includes a set of event message
16 producers, and maintains an event-indication queue of those event messages provided by
17 the event producers using a set of pre-allocated resources. This first aspect allows the
18 system to maintain the event-indication queue even when the event message indicates
19 that allocation of new resources (such as for recording event messages) is unstable. An
20 event-distribution engine distributes event messages to intended recipients and, after
21 having received confirmation that the event messages were received, removes them from
22 the event-indication queue. Recipients (also called "consumers") of event messages re-

1 receive the event messages, acknowledge their receipt thereof, and might take action in re-
2 sponse to the event message.

3
4 In a second aspect of the invention, the system includes a persistent mem-
5 ory (such as NVRAM or other non-volatile memory), in which event messages can be re-
6 corded until they are completely handled by the event-distribution engine. Upon recov-
7 ery from a system crash or other system error, a replay-event producer retrieves those
8 event messages recorded in the persistent memory and not yet completely handled, and
9 re-presents them as event messages for the event-indication queue and the event-
10 distribution engine.

11
12 In a third aspect of the invention, the system includes an initialization
13 memory (also called a "system boot memory"), in which event messages can be recorded
14 until the system has completed its operations of initializing and becoming completely
15 prepared to handle event messages. Upon recovery from a system crash or other system
16 error, the replay-event producer retrieves those event messages recorded in the initializa-
17 tion memory and re-presents them as event messages for the event-indication queue and
18 the event-distribution engine.

19
20 In a fourth aspect of the invention, a cluster of file servers collectively
21 forming a highly-available system shares persistent memories. Each individual file
22 server writes event messages to both its own and at least one other persistent memory, so

1 that upon a system crash or other system error, at least one other file server has a record
2 of those event messages that were presented by event producers but not yet completely
3 handled by the event-distribution engine. Upon indication of a system crash or other
4 system error by a first file server in the cluster, a second file server in the cluster uses its
5 replay-event producer. Thus, the second file server retrieves those event messages re-
6 corded in the persistent memory and not yet completely handled by the first file server,
7 and re-presents them as event messages for its own event-indication queue and event-
8 distribution engine.

9
10 In a fifth aspect of the invention, a cluster of event recipients can be cou-
11 pled to a single multiplexing recipient that includes a second persistent memory in which
12 event messages are recorded after receipt and before being redistributed or otherwise
13 completely handled. Upon recovery from a system crash or other system error by the
14 multiplexing recipient, the multiplexing recipient includes a replay-event producer that
15 retrieves those event messages recorded in the second persistent memory, and re-presents
16 them as event messages as if newly received from an event producer.

17
18 The invention provides an enabling technology for a wide variety of appli-
19 cations for persistent and reliable delivery of event messages, so as to obtain substantial
20 advantages and capabilities that are novel and non-obvious in view of the known art. Ex-
21 amples described below primarily relate to reliable file systems, but the invention is

1 broadly applicable to many different types of systems in which persistent and reliable de-
2 livery of event messages is desired.

3 4 BRIEF DESCRIPTION OF THE DRAWINGS

5
6 Figure 1 shows a block diagram of a portion of a system capable of persis-
7 tent and reliable delivery of event messages.

8
9 Figure 2 shows a process flow diagram of a method for operating a system
10 as in figure 1.

11 12 DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

13
14 In the following description, a preferred embodiment of the invention is de-
15 scribed with regard to preferred process steps and data structures. Embodiments of the
16 invention can be implemented using general-purpose processors or special purpose proc-
17 essors operating under program control, or other circuits, adapted to particular process
18 steps and data structures described herein. Implementation of the process steps and data
19 structures described herein would not require undue experimentation or further invention.

1 *Related Applications*

2
3 Inventions described herein can be used in conjunction with technology de-
4 scribed in the following documents.

- 5
6 *Page 1* • ~~U.S. Patent Application Serial No. _____, Express Mail Mailing No. EL~~
7 ~~524781089US, filed August 18, 2000, in the name of Blake LEWIS, attorney~~
8 ~~docket number 103.1033.01, titled "Reserving File System Blocks"~~
9
10 • U.S. Patent Application Serial No. _____, Express Mail Mailing No.
11 EL524780242US, filed August 18, 2000, in the name of Rajesh SUNDARAM,
12 attorney docket number 103.1034.01, titled "Dynamic Data Storage"
13
14 • U.S. Patent Application Serial No. _____, Express Mail Mailing No.
15 EL524780239US, filed August 18, 2000, in the name of Blake LEWIS, attorney
16 docket number 103.1035.01, titled "Instant Snapshot"
17
18 • U.S. Patent Application Serial No. _____, Express Mail Mailing No.
19 EL524781092US, filed August 18, 2000, in the name of Douglas P. DOUCETTE,
20 attorney docket number 103.1045.01, titled "Improved Space Allocation in a
21 ~~Write Anywhere File System"~~
22

1 and

2 ~~U.S. Patent Application Serial No. _____, Express Mail Mailing No.~~
 3 ~~EL524780256US, filed August 18, 2000, in the name of Ray CHEN, attorney~~
 4 ~~docket number 103.1047.01, titled "manipulation of Zombie Files and Evil-Twin~~
 5 ~~Files"~~
 6

7 Each of these documents is hereby incorporated by reference as if fully set
 8 forth herein. This application claims priority of each of these documents. These docu-
 9 ments are collectively referred to as the "Incorporated Disclosures."

11 *Lexicography*

12
 13 The following terms refer or relate to aspects of the invention as described
 14 below. The descriptions of general meanings of these terms are not intended to be limit-
 15 ing, only illustrative.

- 16
- 17 • *event messages* --- In general, an event message refers to an alert or notification of
 - 18 a system event, for which an intended recipient of that event message may wish to
 - 19 receive and possibly take action in response to. Examples of event messages in-
 - 20 clude notification of errors, or of events to be logged or otherwise administratively
 - 21 monitored.
 - 22

- 1 • *event replay* --- In general, re-presenting event messages (from a memory or other
2 record) by a system element other than the one that generated the original event
3 message. Event messages that are replayed are treated substantially identically to
4 originally generated event messages.
5
- 6 • *persistent, reliable* --- In general, persistent refers to memory or another record
7 that is capable of surviving a disruption such as a system crash or a system error.
8 In general, reliable refers to a process that is capable of being performed atomi-
9 cally or otherwise completely even in the event of a disruption such as a system
10 crash or a system error.
11
- 12 • *pre-allocated resources* --- In general, resources that have been allocated ahead of
13 time for use by a system element (such as for delivery of event messages), so that
14 system element can operate even when disruptions in system operation make it
15 uncertain that resource allocation will be effective at the time when the system
16 element needs those resources.
17
- 18 • *system crashes, system errors* --- In general, a system crash or a system error re-
19 fers to a disruption in system operation sufficiently serious as to place the con-
20 tinuation of system operations (such as delivery of event messages) in doubt. In
21 the description herein, it is assumed that continuation of system operations does
22 not generally survive such disruptions.

- 1
- 2 • *system initialization* --- In general, a time during which a system (such as a file
3 server) is booting or initializing, and so is not necessarily able to perform all nec-
4 essary operations (such as allocation of resources or delivery of event messages).
- 5

6 As noted above, these descriptions of general meanings of these terms are
7 not intended to be limiting, only illustrative. Other and further applications of the inven-
8 tion, including extensions of these terms and concepts, would be clear to those of ordi-
9 nary skill in the art after perusing this application. These other and further applications
10 are part of the scope and spirit of the invention, and would be clear to those of ordinary
11 skill in the art, without further invention or undue experimentation.

12

13 *System Elements*

14

15 Figure 1 shows a block diagram of a portion of a system capable of persis-
16 tent and reliable delivery of event messages.

17

18 A system 100 includes a set of event producers 110, a set of pre-allocated
19 initialization event message resources 120, a set of pre-allocated post-initialization event
20 message resources 130, a persistent memory 140, an event indication queue 150, an event
21 distribution engine 160, a set of event recipients 170 including at least one multiplexing

1 recipient 171 and a set of intended recipients 172, a second persistent memory 180 at the
2 multiplexing recipient 171, and an event replay engine 190.

3
4 The event producers 110 include system elements, such as software mod-
5 ules or hardware circuits, each capable of generating at least one event message 111 for
6 delivery to at least one intended recipient 172. In a preferred embodiment, each event
7 message 111 has a standardized format, including information about the time the event
8 was recognized, the system element that recognized the event, the nature of the event,
9 and any detailed information about the event necessary or desirable for the intended re-
10 cipient 172 to know.

11
12 The pre-allocated initialization event message resources 120 include mem-
13 ory, and possibly other resources, for recording and maintaining event messages for fur-
14 ther processing. In a preferred embodiment, the pre-allocated initialization event mes-
15 sage resources 120 include resources allocated by the system 100 prior to generation of
16 any new event messages 111. In a preferred embodiment, the event messages 111 are re-
17 corded in the pre-allocated initialization event message resources 120 before being proc-
18 essed by the event distribution engine 160.

19
20 The pre-allocated post-initialization event message resources 130, similar
21 to the pre-allocated initialization event message resources 120, memory, and possibly
22 other resources, for recording and maintaining event messages for further processing. In

1 a preferred embodiment, the pre-allocated post-initialization event message resources 130
2 include resources allocated by the system 100 during an initialization period and prior to
3 generation of any new event messages 111. In a preferred embodiment, event messages
4 111 that are generated during the system initialization period are recorded in the pre-
5 allocated post-initialization event message resources 130. After the system initialization
6 period, event messages 111 recorded therein are processed by the event distribution en-
7 gine 160.

8
9 The persistent memory 140 includes a memory, such as NVRAM, SRAM,
10 or other memory whose contents are expected to survive a system crash or system error.
11 In a preferred embodiment, the persistent memory 140 includes NVRAM also used with
12 the WAFL file system. However, in alternative embodiments, the persistent memory 140
13 may include any other form of persistent memory, whether NVRAM or not, and whether
14 coordinating with aspects of the WAFL file system or not. Thus, upon recovery from a
15 system crash or system error, the persistent memory 140 will still record those event mes-
16 sages 111 that were not fully processed before the system crash or system error.

17
18 The event indication queue 150 includes a memory having a queue of in-
19 formation about event messages 111 (such as the event messages 111 themselves).

20
21 The event distribution engine 160 includes a system element capable of
22 reading information about event messages 111 from the event indication queue 150 and

1 capable of delivering at those event messages 111 to intended recipients 172 thereof. In a
2 preferred embodiment, the event distribution engine 160 includes a software module in
3 the system 100.

4
5 The event recipients 170 (including at least one multiplexing recipient 171
6 and a set of intended recipients 172) include system elements, possibly at remote devices
7 such as clients for the file server system 100, capable of receiving event messages 111
8 and deciding whether or not to act in response to those event messages 111. In a pre-
9 ferred embodiment, actions take with regard to event messages 111 can include alerts or
10 notification of selected users (such as a system operator), logging the event messages
11 111, or maintaining statistics with regard thereto.

12
13 The second persistent memory 180 at the multiplexing recipient 171 in-
14 cludes, similar to the persistent memory 140, a memory, such as NVRAM, SRAM, or
15 other memory whose contents are expected to survive a system crash or system error.
16 The multiplexing recipient 171 includes a recipient replay element 181, capable of read-
17 ing information about event messages 111 from the second persistent memory 180 and
18 capable of replaying those event messages 111 as if newly received by the multiplexing
19 recipient 171 (thus delivering those event messages 111 to the intended recipients 172).

20
21 The event replay engine 190 includes a system element capable of reading
22 information about event messages 111 from the persistent memory 140 and capable of

1 replaying those event messages 111 as if newly generated. The replay element 190 in-
2 cludes a system initialization replay sub-element 191, an incomplete event distribution
3 replay sub-element 192, and a cooperating systems replay sub-element 193.

4 5 *Method of Operation*

6
7 Figure 2 shows a process flow diagram of a method for operating a system
8 as in figure 1.

9
10 A method 200 includes a set of flow points and a set of steps. The system
11 100 performs the method 200. Although the method 200 is described serially, the steps of
12 the method 200 can be performed by separate elements in conjunction or in parallel,
13 whether asynchronously, in a pipelined manner, or otherwise. There is no particular re-
14 quirement that the method 200 be performed in the same order in which this description
15 lists the steps, except where so indicated.

16
17 As described below, the method 200 includes a set of processes, each of
18 which has a set of tasks operating independently and asynchronously with regard to each
19 other.

20 21 1. *Processing Event Messages*

22

1 A first process in the method 200 is described with regard to a flow point
2 210, a flow point 220, and steps there-between. This first process includes a set of three
3 tasks, each of which operates independently and asynchronously with regard to each
4 other.

5
6 At the flow point 210, the system 100 is ready to receive an event message
7 111.

8 9 *Event Generation*

10
11 A first task includes a sequence including a step 211, a step 212, a step 213,
12 and a step 214. In a preferred embodiment, the first task in its process includes these
13 steps being performed in sequence and repeatedly.

14
15 At the step 211, an event producer 110 generates an event message 111.

16
17 At the step 212, if the system 100 is in normal operation, the event message
18 111 is recorded in the pre-allocated initialization event message resources 120. If the
19 system 100 is still in its initialization time duration, the event message 111 is recorded in
20 the pre-allocated post-initialization event message resources 130.

1 At the step 213, the system 100 copies information about the event message
2 111 to a set of locations in the persistent memory 140. In a preferred embodiment, the
3 persistent memory 140 includes a set of memory sections 141, each persistently main-
4 taining system information; at least one of these memory sections 141 maintains infor-
5 mation about event messages 111. Other memory sections 141 maintain information
6 about other aspects of the system 100, such as a consistency state of the file system or a
7 set of incomplete file system requests.

8
9 In a preferred embodiment, the memory sections 141 associated with event
10 messages 111 maintain information about those event messages 111 in a FIFO having a
11 head pointer and a tail pointer. FIFOs are known in the art of computer data storage.
12 When information about a new event message 111 is recorded in the persistent memory
13 140, the FIFO is updated to add the information about the new event message 111 to an
14 end of the list. When confirmation is received that the event message 111 was delivered
15 to its intended recipients 172, the FIFO is updated to remove the information about the
16 event message 111.

17
18 At the step 214, similar to the step 213, the system 100 copies information
19 about the event message 111 to the event indication queue 150. In a preferred embodi-
20 ment, the event indication queue 150 includes a FIFO similar to that maintained in the
21 persistent memory 140.

Event Distribution

A second task includes a sequence including a step 215. In a preferred embodiment, the second task in its process includes this step being performed repeatedly.

At the step 215, the event distribution engine 160 responds to the information about the event message 111 in the event indication queue 150. The event distribution engine 160 delivers the event message 111 to its intended recipients 172. As part of this step, each particular intended recipient 172, when it receives the event message 111, responds to the event distribution engine 160 to confirm its receipt of the event message 111.

Event Confirmation

A third task includes a sequence including a step 216 and a step 217. In a preferred embodiment, the third task in its process includes these steps being performed in sequence and repeatedly.

At the step 216, the event distribution engine 160 awaits confirmation from each intended recipient 172 that the event message 111 was received by that particular intended recipient 172. When the event distribution engine 160 receives confirmation from all intended recipients 172, the method proceeds with the next step.

At the step 217, the event distribution engine 160 removes the information about the event message 111 from the event indication queue 150 and from the persistent memory 140.

At the flow point 220, the system 100 has completely processed the event message 111.

2. *Replaying Event Messages*

At a flow point 230, the system 100 has recovered from a system crash or a system error.

At a step 231, the event replay engine 190 reads information about event messages 111 from the persistent memory 140. As part of this step, the event replay engine 190 performs three sub-steps 231(a), 231(b), and 231(c).

At the sub-step 231(a), the system initialization replay sub-element 191 reads information about event messages 111 associated with the pre-allocated post-initialization event message resources 130. The event replay engine 190 replays these event messages 111.

1 At the sub-step 231(b), the incomplete event distribution replay sub-
2 element 192 reads information about event messages 111 associated with the pre-
3 allocated initialization event message resources 120. The event replay engine 190 re-
4 plays these event messages 111 next.

5
6 At the sub-step 231(c), the cooperating systems replay sub-element 193
7 reads information about event messages 111, from the persistent memory 140, associated
8 with and stored there by a cooperating system 100. The event replay engine 190 replays
9 these event messages 111 only if the cooperating system 100 is not operational at the
10 time.

11
12 In a preferred embodiment, multiple cooperating systems 100 (preferably a
13 pair of exactly two) are each capable of reading and writing to each other's persistent
14 memories 140. Thus, when a first cooperating system 100 in the pair writes to its persis-
15 tent memory 140, the second cooperating system 100 in the pair is able to read from that
16 persistent memory 140. If the first cooperating system 100 suffers a system crash or
17 system error, the second cooperating system 100, upon recognizing that system crash or
18 system error, proceeds to replay the event messages 111 from the first cooperating sys-
19 tem's persistent memory 140. Operation of multiple cooperating systems 100 is further
20 described in the Incorporated Disclosures, particularly with regard to techniques used to
21 prevent multiple cooperating systems 100 from disrupting each other's operation.

1 As noted herein, “replay” of event messages 111 is treated by the event in-
2 dication queue 150 and the event distribution engine 160 as if the event messages 111
3 were newly generated. Replayed event messages 111 are processed and delivered before
4 new event messages 111, according to the portion of the method 200 described with re-
5 gard to flow point 210 and flow point 220.

6
7 At a flow point 240, the system 100 has replayed all event messages 111
8 not yet fully processed, and is ready to proceed at the flow point 210.

9
10 3. *Multiplexing Recipient Operation*

11
12 A third process in the method 200 is described with regard to a flow point
13 250, a flow point 260, and steps there-between. Similar to the first process, this third
14 process includes a set of three tasks, each of which operates independently and asynchro-
15 nously with regard to each other.

16
17 At the flow point 250, a multiplexing recipient 171 is ready to receive an
18 event message 111.

19
20 *Event Reception*

1 A first task includes a sequence including a step 251, a step 252, and a step
2 253. In a preferred embodiment, this first task in its process includes these steps being
3 performed in sequence and repeatedly.

4
5 At the step 251, the multiplexing recipient 171 receives the event message
6 111 from the event distribution engine 160.

7
8 At the step 252, similarly to the steps described with regard to the flow
9 point 210 and the flow point 220, the multiplexing recipient 171 records information
10 about the event message 111 in its second persistent memory 180.

11
12 At the step 253, the multiplexing recipient 171 (optionally) responds to the
13 event message 111 by confirming that it was received at the multiplexing recipient 171
14 (but not necessarily at the intended recipients 172).

15
16 *Event Multiplexing*

17
18 A second task includes a sequence including a step 254. In a preferred em-
19 bodiment, this second task in its process includes this step being performed repeatedly.

20
21 At the step 254, the multiplexing recipient 171 (optionally) determines to
22 which intended recipients 172 to deliver the event message 111. In a preferred embodi-

ment, the multiplexing recipient 171 filters the event messages 111 it receives, so that it delivers only those event messages 111 it receives to their actual intended recipients 172. For example, a particular intended recipient 172 might determine that it is only interested in a particular subclass of event messages 111. In such cases, the multiplexing recipient 171 delivers only that particular subclass of event messages 111 to that particular intended recipient 172.

Event Confirmation

A third task includes a sequence including a step 255 and a step 256. In a preferred embodiment, this third task in its process includes these steps being performed in sequence and repeatedly.

At the step 255, the multiplexing recipient 171 awaits confirmation from each particular intended recipient 172 that the particular intended recipient 172 has received the event message 111.

At the step 256, the multiplexing recipient 171 receives such confirmation from individual intended recipients 172. As part of this step, the multiplexing recipient 171 (optionally) forwards those confirmations on to the event distribution engine 160. When the multiplexing recipient 171 receives all such confirmations, it removes the information about the event message 111 from the second persistent memory 180.

At the flow point 260, the multiplexing recipient 171 has completely processed the event message 111.

4. *Replaying Multiplexed Event Messages*

At a flow point 270, the multiplexing recipient 171 has recovered from a system crash or a system error.

At a step 271, the multiplexing recipient 171 reads information about event messages 111 from the second persistent memory 180.

At a step 272, the multiplexing recipient 171 replays the event messages 111 from the second persistent memory 180.

“Replay” of event messages 111 by the multiplexing recipient 171 is similar to replay of event messages as described above with regard to the event indication queue 150 and the event distribution engine 160.

At a flow point 280, the multiplexing recipient 171 has replayed all event messages 111 not yet fully processed, and is ready to proceed at the flow point 250.

1 5. *Confirming Event Messages*

2
3 As described above, there are steps at which the system 100 or the multi-
4 plexing recipient 171 awaits confirmation of the event message 111 from the intended re-
5 cipient 172. In a preferred embodiment, confirmation of event messages 111 is per-
6 formed by each intended recipient 172 as described below with regard to a flow point
7 290, a flow point 300, and steps there-between.

8
9 At the flow point 290, the intended recipient 172 is ready to receive an
10 event message 111.

11
12 At a step 291, the intended recipient 172 receives an event message 111.

13
14 At a step 292, the intended recipient 172 parses the event message 111 and
15 processes the event message 111 according to its own (internal) processing rules for that
16 event message 111.

17
18 At a step 293, the intended recipient 172 generates a confirmation message
19 and sends that confirmation message to the sender of the event message 111.

20
21 At the flow point 300, the intended recipient 172 has received, processed,
22 and confirmed the event message 111. The sender of the event message 111, upon re-

1 ceipt of the confirmation message, can regard the event message 111 as completely han-
2 dled and can safely delete it.

3
4 *Generality of the Invention*

5
6 The invention has general applicability to various fields of use, not neces-
7 sarily related to the services described above. For example, these fields of use can include
8 one or more of, or some combination of, the following:

- 9
- 10 • The invention is applicable to persistent and reliable delivery of messages other
11 than event messages.
 - 12
 - 13 • The invention is applicable to persistent and reliable operation of other processes
14 than delivery of messages.
 - 15
 - 16 • The invention is applicable to mutually cooperating systems to perform other per-
17 sistent and reliable operations.
 - 18
 - 19 • The invention is applicable to hierarchical cooperating systems to perform other
20 persistent and reliable operations.
 - 21

